

## Topic 6: - Topic Analysis

Clarissa Boyajian

2022-05-09

```
## -- read in, clean, and wrangle data -- ##
comments_df<-read_csv(here("data/comments_df.csv"))

epa_corp <- corpus(x = comments_df, text_field = "text")
epa_corp.stats <- summary(epa_corp)
tokens <- tokens(epa_corp, remove_punct = TRUE, remove_numbers = TRUE)

# I added some project-specific stop words here
add_stops <- c(stopwords("en"),"environmental", "justice", "ej", "epa", "public", "comment")
tokens_1 <- tokens_select(tokens, pattern = add_stops, selection = "remove")

doc_feature_matrix_common <- dfm(tokens_1, tolower = TRUE)
doc_feature_matrix <- dfm_wordstem(doc_feature_matrix_common)
doc_feature_matrix <- dfm_trim(doc_feature_matrix, min_docfreq = 2) #remove terms only appearing in one doc

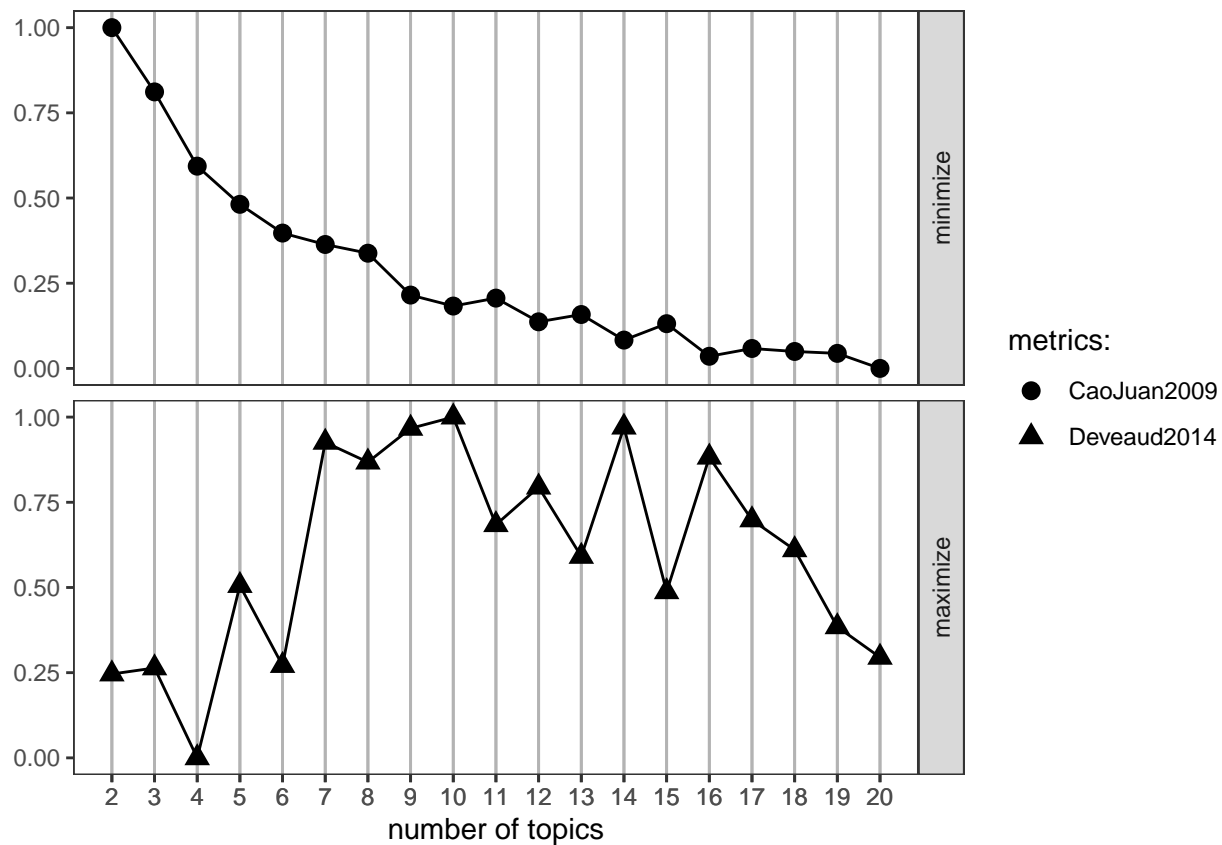
# remove rows (docs) with all zeros
select_idx <- slam::row_sums(doc_feature_matrix) > 0
doc_feature_matrix <- doc_feature_matrix[select_idx, ]
```

Calculate value of k that is most likely

```
# calculate what initial value of k is most likely
result <- FindTopicsNumber(
  doc_feature_matrix,
  topics = seq(from = 2, to = 20, by = 1),
  metrics = c("CaoJuan2009", "Deveaud2014"),
  method = "Gibbs",
  control = list(seed = 77),
  verbose = TRUE
)
```

```
## fit models... done.
## calculate metrics:
##   CaoJuan2009... done.
##   Deveaud2014... done.
```

```
FindTopicsNumber_plot(result)
```



## Model 1 (k = 10)

```
k <- 10

topicModel_k10 <- LDA(doc_feature_matric, k,
  method = "Gibbs",
  control = list(iter = 500, verbose = 25))

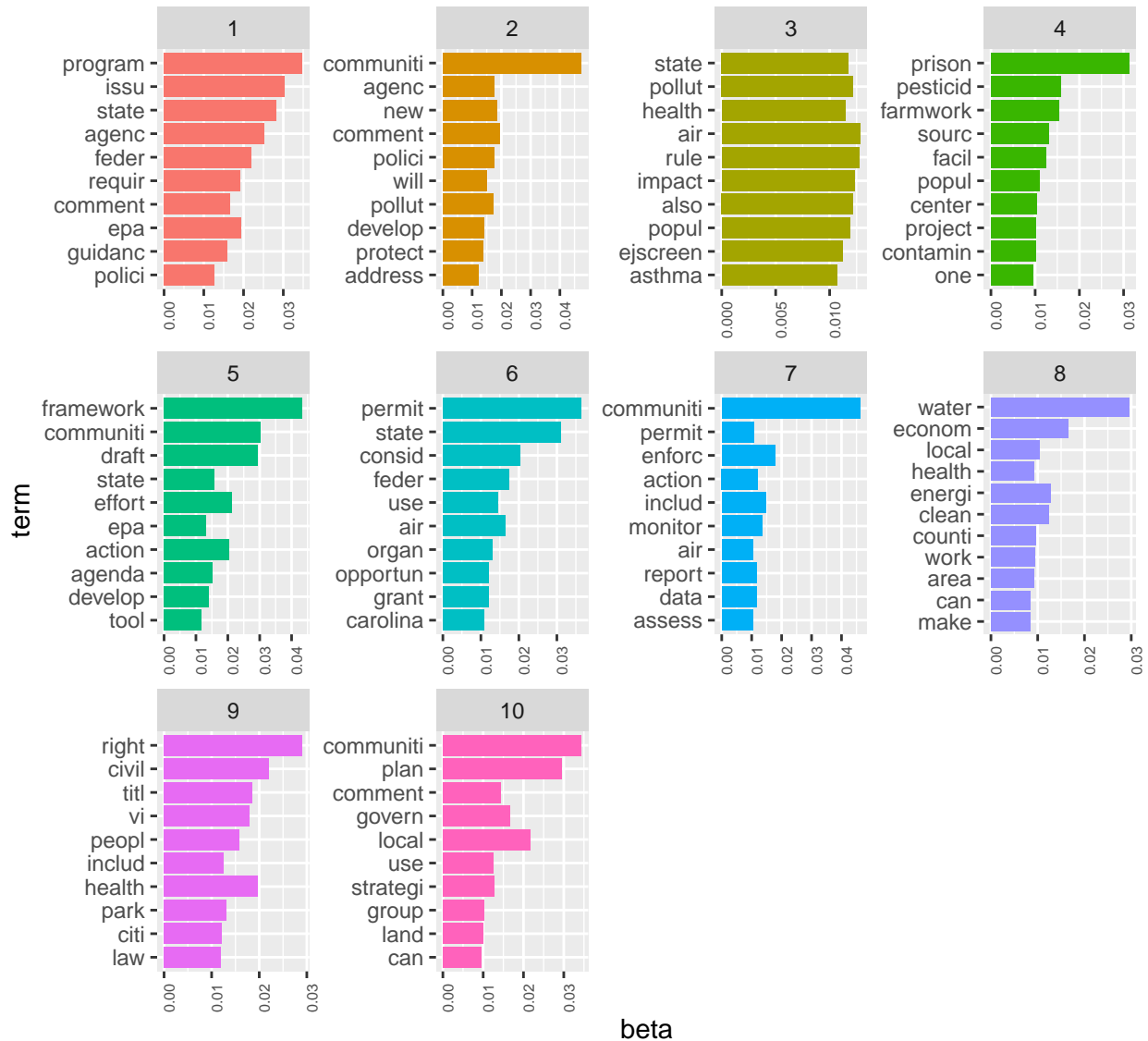
tmResult_k10 <- posterior(topicModel_k10)
theta_k10 <- tmResult_k10$topics
beta_k10 <- tmResult_k10$terms # probability of each term in each topic
vocab_k10 <- (colnames(beta_k10))

comment_topic_k10 <- tidy(topicModel_k10, matrix = "beta")

top_terms_k10 <- comment_topic_k10 %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
  arrange(topic, -beta)

top_terms_k10 %>%
  mutate(term = reorder(term, beta)) %>%
  ggplot(aes(term, beta, fill = factor(topic))) +
```

```
geom_col(show.legend = FALSE) +
facet_wrap(~ topic, scales = "free") +
coord_flip() +
theme(axis.text.x = element_text(angle = 90, hjust = 1, size = 6))
```



## Model 2 (k = 16)

```
k <- 16

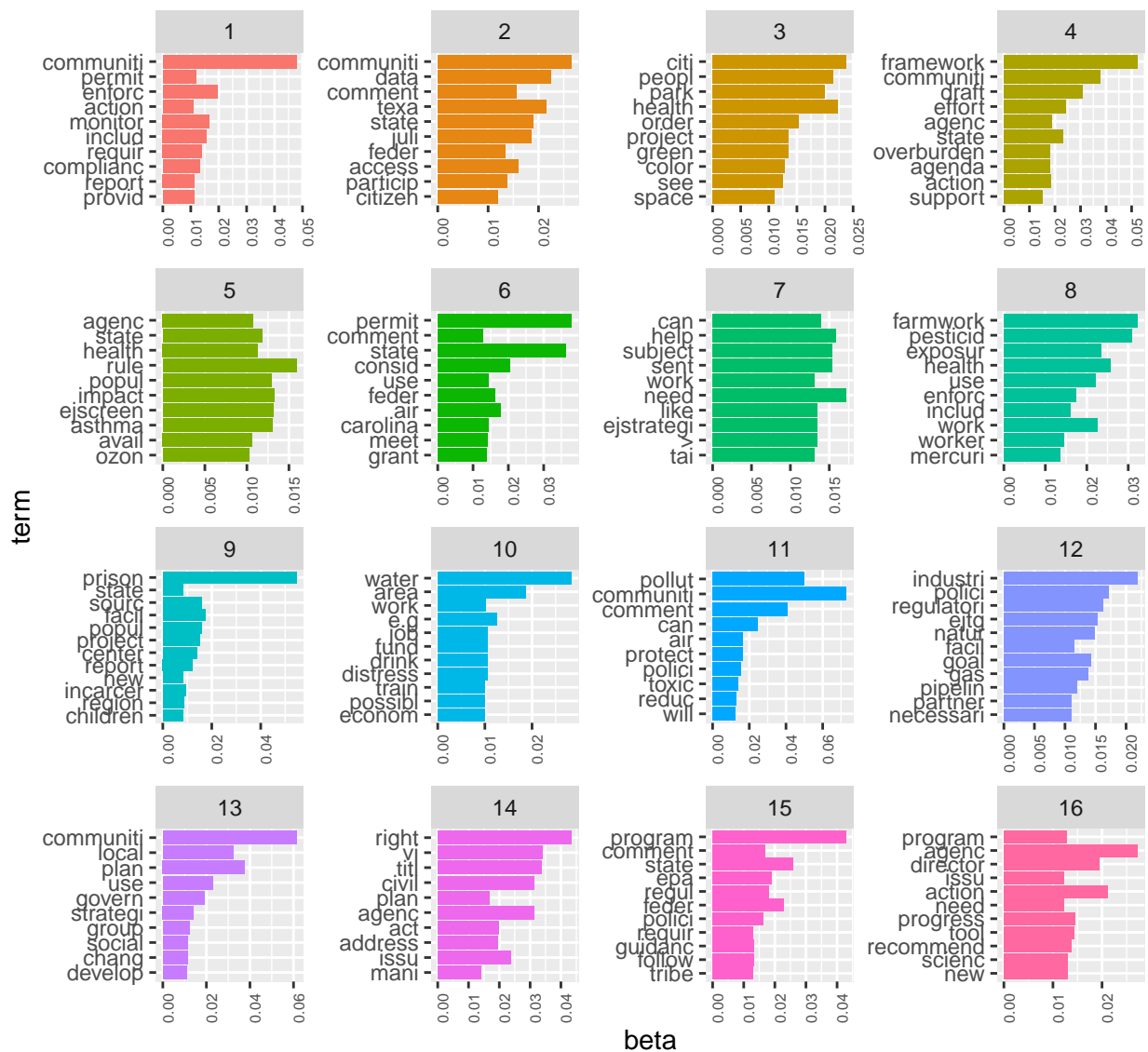
topicModel_k16 <- LDA(doc_feature_matric, k,
                      method = "Gibbs",
                      control = list(iter = 500, verbose = 25))

tmResult_k16 <- posterior(topicModel_k16)
theta_k16 <- tmResult_k16$topics
beta_k16 <- tmResult_k16$terms # probability of each term in each topic
vocab_k16 <- (colnames(beta_k16))

comment_topic_k16 <- tidy(topicModel_k16, matrix = "beta")

top_terms_k16 <- comment_topic_k16 %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
  arrange(topic, -beta)

top_terms_k16 %>%
  mutate(term = reorder(term, beta)) %>%
  ggplot(aes(term, beta, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  coord_flip() +
  theme(axis.text.x = element_text(angle = 90, hjust = 1, size = 6))
```



### Model 3 (k = 20)

```
k <- 20

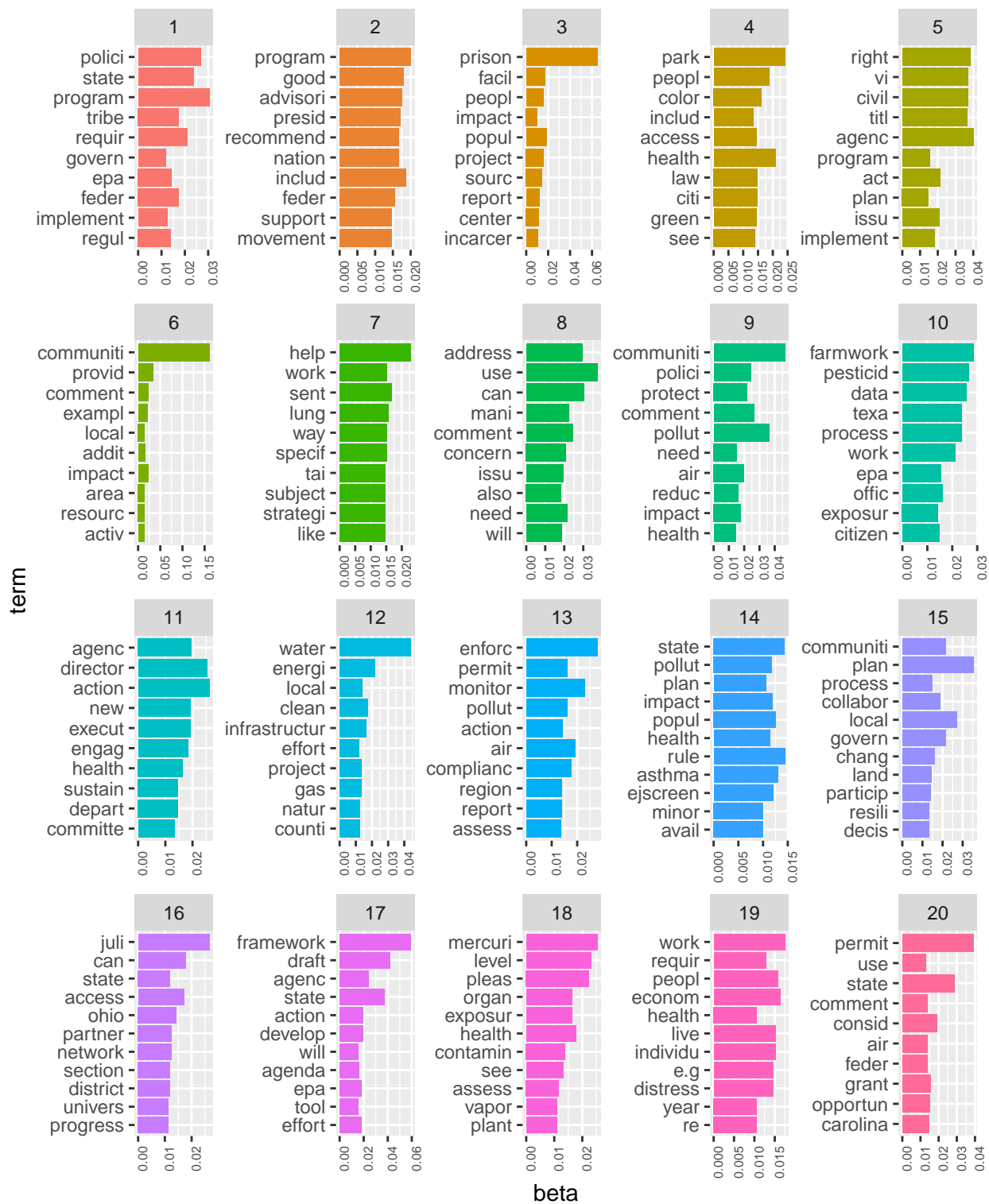
topicModel_k20 <- LDA(doc_feature_matric, k,
                      method = "Gibbs",
                      control = list(iter = 500, verbose = 25))

tmResult_k20 <- posterior(topicModel_k20)
theta_k20 <- tmResult_k20$topics
beta_k20 <- tmResult_k20$terms # probability of each term in each topic
vocab_k20 <- (colnames(beta_k20))

comment_topic_k20 <- tidy(topicModel_k20, matrix = "beta")

top_terms_k20 <- comment_topic_k20 %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
  arrange(topic, -beta)

top_terms_k20 %>%
  mutate(term = reorder(term, beta)) %>%
  ggplot(aes(term, beta, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  coord_flip() +
  theme(axis.text.x = element_text(angle = 90, hjust = 1, size = 6))
```



## Pick best k value

```
docIds <- c(1, 2, 3, 4, 5, 6)
N <- length(docIds)

## -- k = 10 -- ##
# name each topic based on 1st 5 terms in topic
top5termsPerTopic_k10 <- terms(topicModel_k10, 5)
topicNames_k10 <- apply(top5termsPerTopic_k10, 2, paste, collapse = " ")

# get topic proportions from example documents
topicProportions_K10 <- theta_k10[docIds,]
colnames(topicProportions_K10) <- topicNames_k10

vizDataFrame_k10 <- melt(cbind(data.frame(topicProportions_K10),
                                document = factor(1:N)),
                        variable.name = "topic",
                        id.vars = "document")

plot_k10 <- ggplot(data = vizDataFrame_k10,
                  aes(x = topic,
                      y = value,
                      fill = document),
                  ylab = "proportion") +
  geom_bar(stat="identity") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1),
        legend.position = "none") +
  coord_flip() +
  facet_wrap(~ document, ncol = N) +
  labs(title = "K=10")

## -- k = 16 -- ##
# name each topic based on 1st 5 terms in topic
top5termsPerTopic_k16 <- terms(topicModel_k16, 5)
topicNames_k16 <- apply(top5termsPerTopic_k16, 2, paste, collapse = " ")

# get topic proportions from example documents
topicProportions_k16 <- theta_k16[docIds,]
colnames(topicProportions_k16) <- topicNames_k16

vizDataFrame_k16 <- melt(cbind(data.frame(topicProportions_k16),
                                document = factor(1:N)),
                        variable.name = "topic",
                        id.vars = "document")

plot_k16 <- ggplot(data = vizDataFrame_k16,
                  aes(x = topic,
                      y = value,
                      fill = document),
                  ylab = "proportion") +
  geom_bar(stat="identity") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1),
        legend.position = "none") +
```



```
coord_flip() +
facet_wrap(~ document, ncol = N) +
labs(title = "K=16")
```

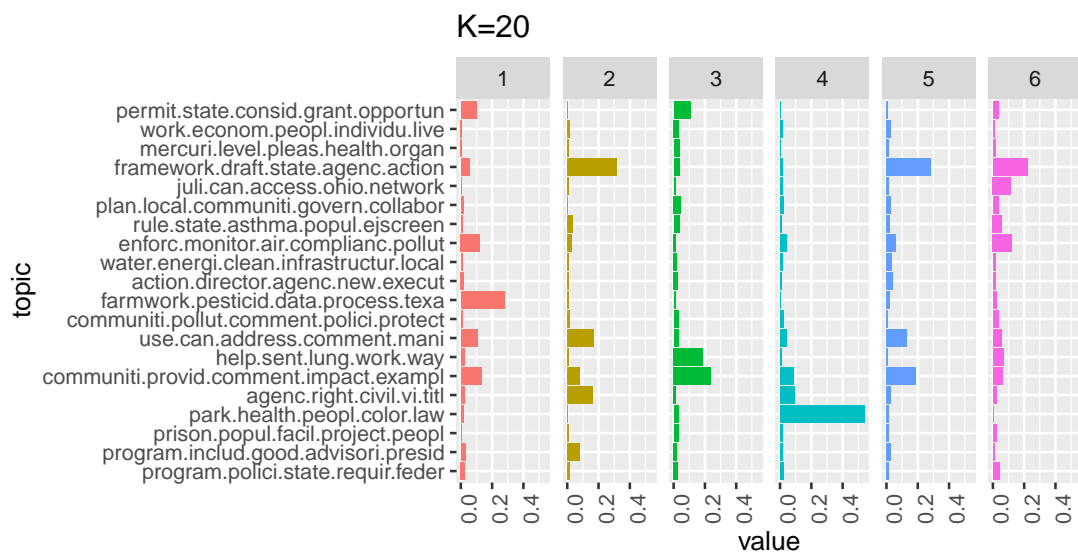
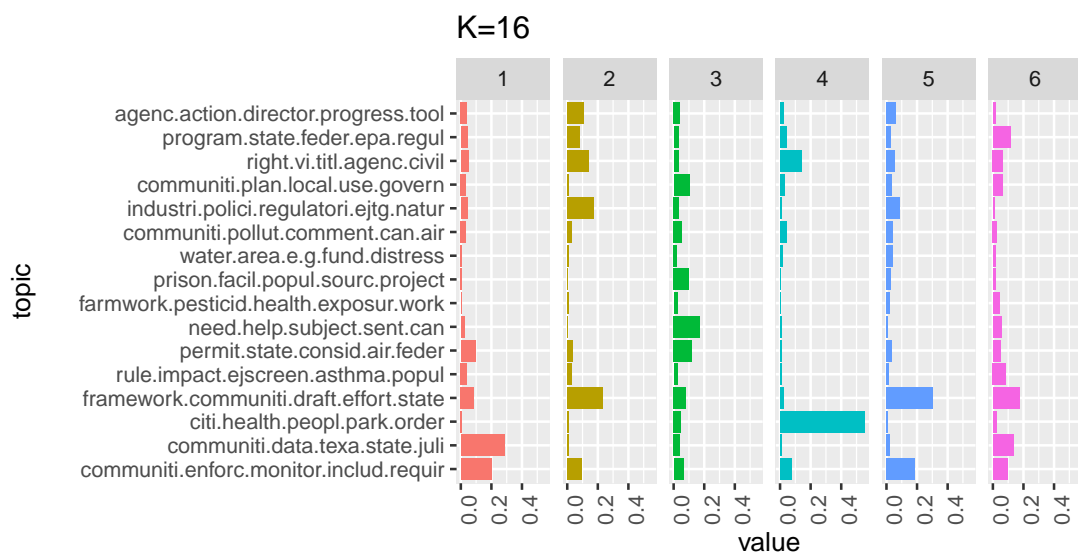
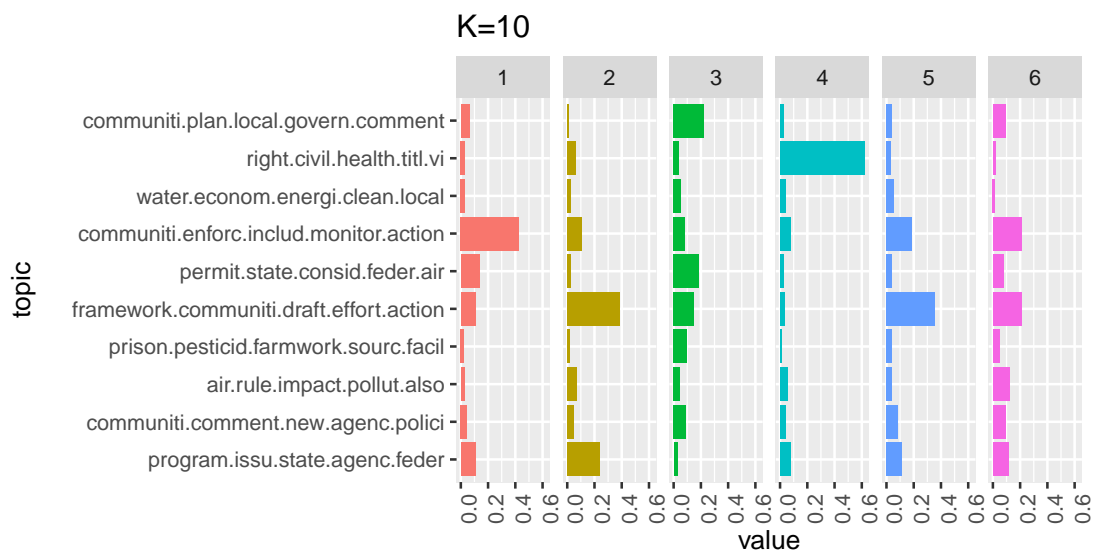
```
## -- k = 20 -- ##
# name each topic based on 1st 5 terms in topic
top5termsPerTopic_k20 <- terms(topicModel_k20, 5)
topicNames_k20 <- apply(top5termsPerTopic_k20, 2, paste, collapse = " ")

# get topic proportions from example documents
topicProportions_k20 <- theta_k20[docIds,]
colnames(topicProportions_k20) <- topicNames_k20

vizDataFrame_k20 <- melt(cbind(data.frame(topicProportions_k20),
                                   document = factor(1:N)),
                        variable.name = "topic",
                        id.vars = "document")

plot_k20 <- ggplot(data = vizDataFrame_k20,
                  aes(x = topic,
                      y = value,
                      fill = document),
                  ylab = "proportion") +
  geom_bar(stat="identity") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1),
        legend.position = "none") +
  coord_flip() +
  facet_wrap(~ document, ncol = N) +
  labs(title = "K=20")
```

```
plot_k10 / plot_k16 / plot_k20
```



```
svd_tsne <- function(x){
  tsne(svd(x)$u)
}
```

```
json_k10 <- createJSON(
  phi = tmResult_k10$terms,
  theta = tmResult_k10$topics,
  doc.length = rowSums(doc_feature_matric),
  vocab = colnames(doc_feature_matric),
  term.frequency = colSums(doc_feature_matric),
  mds.method = svd_tsne,
  plot.opts = list(xlab = "", ylab = ""))
```

```
json_k16 <- createJSON(
  phi = tmResult_k16$terms,
  theta = tmResult_k16$topics,
  doc.length = rowSums(doc_feature_matric),
  vocab = colnames(doc_feature_matric),
  term.frequency = colSums(doc_feature_matric),
  mds.method = svd_tsne,
  plot.opts = list(xlab = "", ylab = ""))
```

```
json_k20 <- createJSON(
  phi = tmResult_k20$terms,
  theta = tmResult_k20$topics,
  doc.length = rowSums(doc_feature_matric),
  vocab = colnames(doc_feature_matric),
  term.frequency = colSums(doc_feature_matric),
  mds.method = svd_tsne,
  plot.opts = list(xlab = "", ylab = ""))
```

```
serVis(json_k20)
```

Based on the above minimization and maximization plots, 10 was the plot with the highest maximization value. We know there are 9 priority areas for the EPA. So it would make sense that there are topics for each priority area, plus one extra topic that catches on additional miscellaneous topic that appears in all PDFs.